

# An Empirical Study of Translation Hypothesis Ensembling with Large Language Models

António Farinhas<sup>1,2</sup>, José G. C. de Souza<sup>3</sup>, André F. T. Martins<sup>1,2,3</sup>

<sup>1</sup>Instituto Superior Técnico (Lisbon ELLIS Unit), <sup>2</sup>Instituto de Telecomunicações, <sup>3</sup>Unbabel

There's lots of research on task-specific NMT models but LLMs offer a new perspective!

We generate multiple hypotheses by using a single prompt and sampling multiple times; we ensemble them using different techniques.

## Translation Hypothesis Ensembling

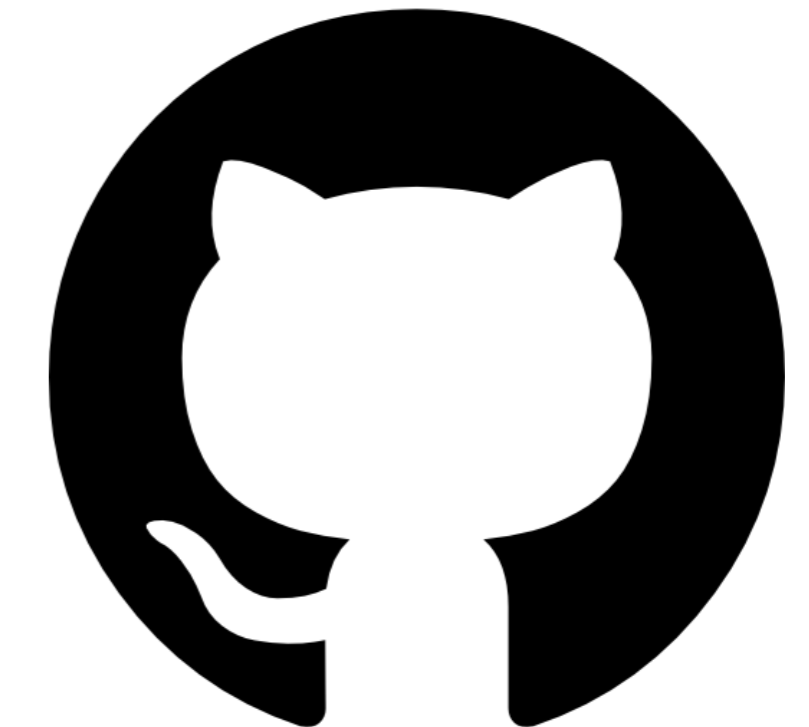
Using external quality estimation/evaluation models

- ranking with QE:  $\hat{y}_{\text{ranking}} = \operatorname{argmax}_{y \in \bar{y}} \operatorname{CometKiwi}(y)$
- MBR decoding:  $\hat{y}_{\text{mbr}} = \operatorname{argmax}_{y \in \bar{y}} \mathbb{E}_{Y \sim p_\theta} [\operatorname{COMET}(Y, y)]$

Using the LLM

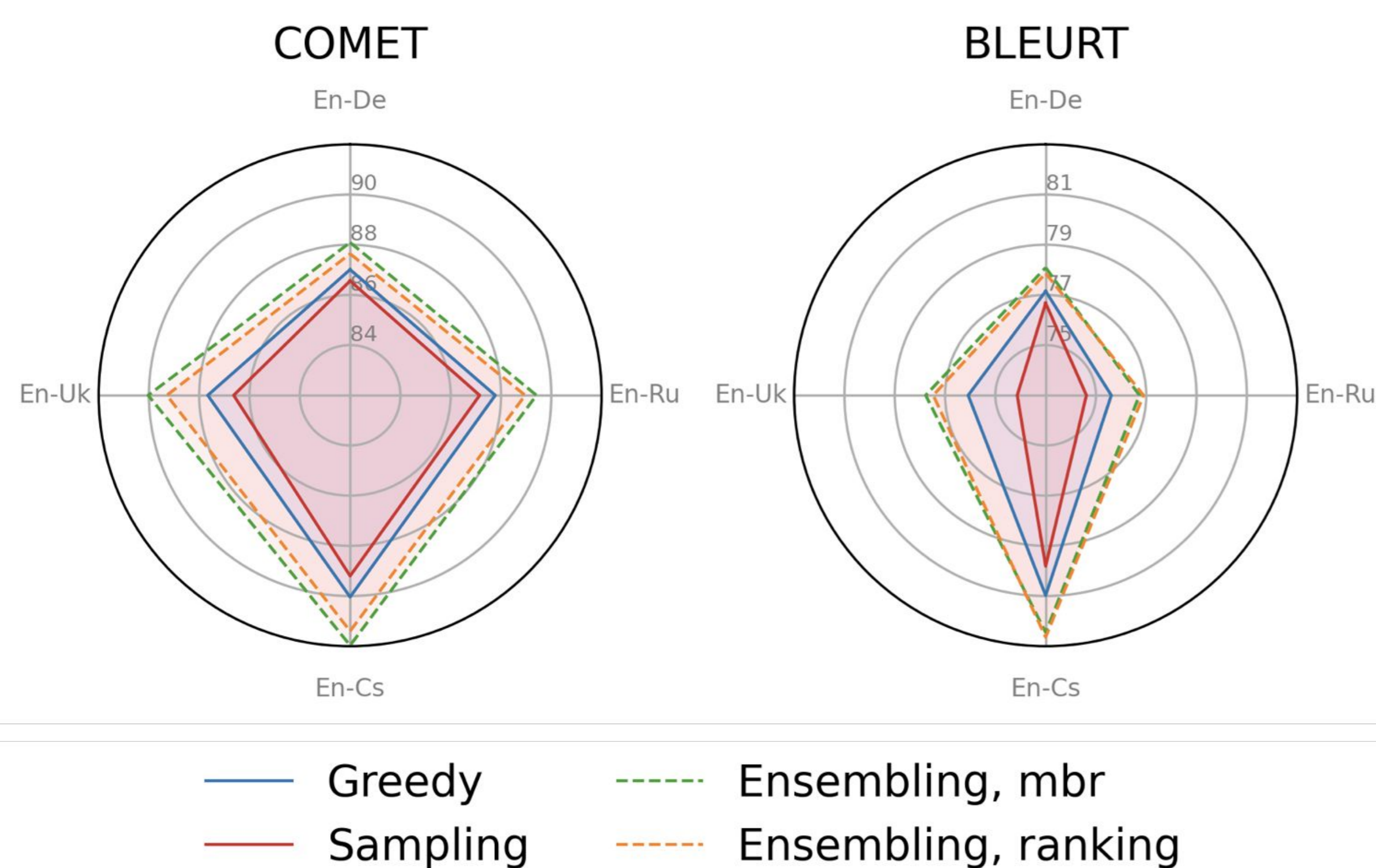
- ChooseBest: formulated as a multiple choice question
- GenerateBest: asking the LLM to generate a final prediction

Check the paper for more LPs and analysis!



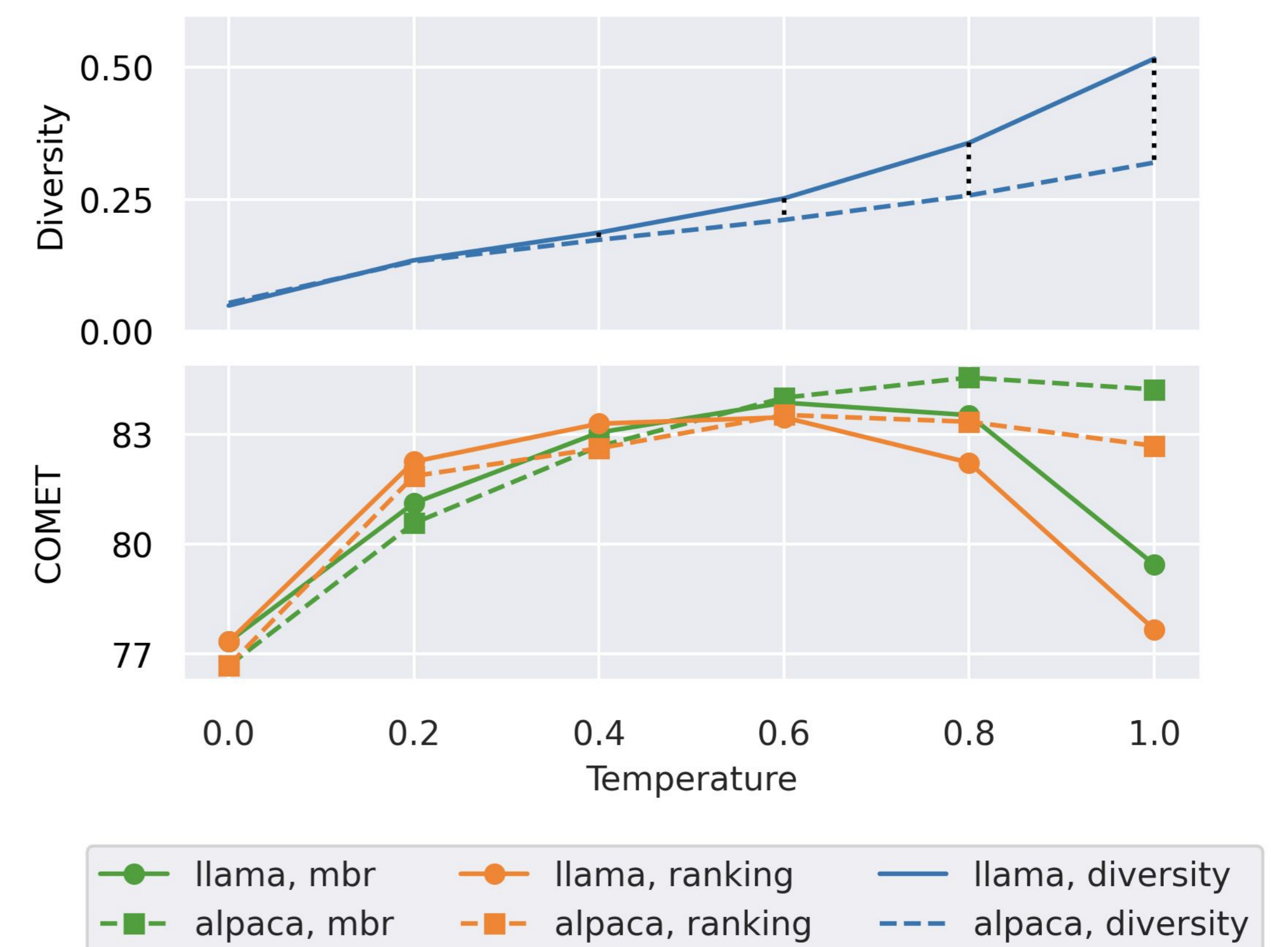
## ChatGPT

- translation quality can be enhanced with a small number of unbiased samples, especially for EN-X



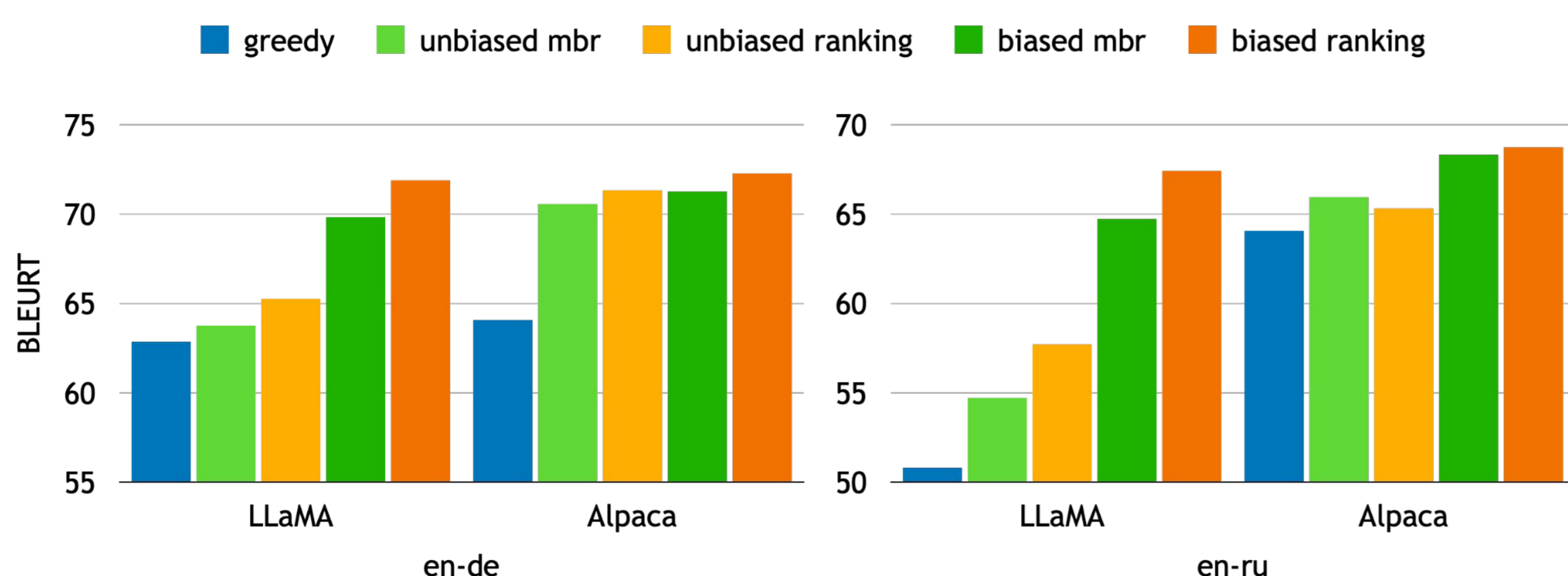
## Biasedness, diversity, and quality

- the diversity between hypotheses increases with the sampling temperature at a different rate for LLaMA and Alpaca



## LLaMA and Alpaca

- ensembles of unbiased samples from LLaMA don't perform well
- alpaca performs better and biasing samples boosts performance



## Hallucinations

- hallucination rate decreases with instruction tuning
- ensembling translations decreases the number of hallucinations

